

CS6501: Deep Learning for Visual Recognition

Seq2Seq Model & Text-to-Image Synthesis



Presenter: Fuwen Tan

Today's Class

- Mini-batch training of the RNN model
 - Special “End-of-Sequence” token: <end>
 - Padding
- Sequence-to-sequence model
 - Neural Machine Translation^[1]
- Text-to-Image Synthesis^[2]

[1] Effective Approaches to Attention-based Neural Machine Translation. Thang Luong, Hieu Pham, and Christopher D. Manning. EMNLP 2015

[2] Text2Scene: Generating Compositional Scenes from Textual Descriptions. Fuwen Tan, Song Feng, Vicente Ordonez. CVPR 2019.

A RNN model will never end

“Hello”, “world”, “!”, “!”, “!”, “!”, “!”, ...

Unless: set the maximum length before hand

Sample 1	<i>"hello"</i>	<i>"world"</i>	<i>"java"</i>	<i>"is"</i>	<i>"better"</i>
Sample 2	<i>"hello"</i>	<i>"hoos"</i>	<i>"I"</i>	<i>"like"</i>	<i>"python"</i>
Sample 3	<i>"one"</i>	<i>"plus"</i>	<i>"eight"</i>	<i>"equals"</i>	<i>"to"</i>

I want sentences of 5 words

Or: learn to predict the END.

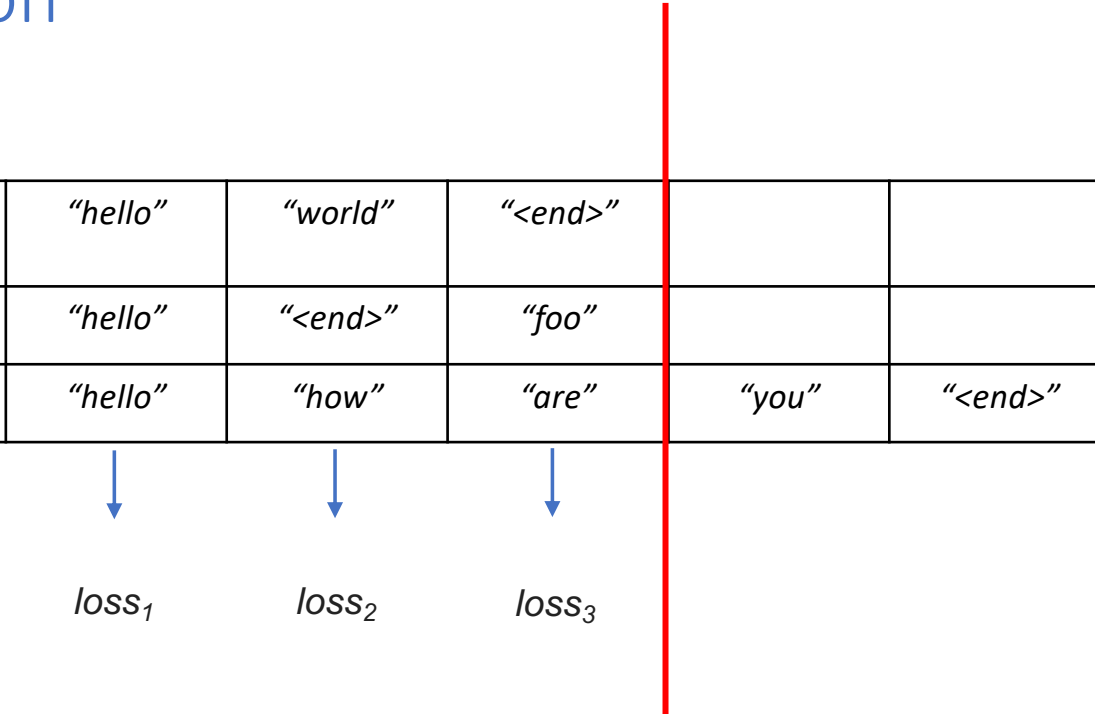
Ground-truth: *“Hello”, “world”, “<end>”*

Training: learn to generate the ground-truth sequence with *“<end>”*.

Testing: generate the sequence until an *“<end>”* is predicted.

Computing loss: what if #ground-truth \neq #prediction

Ground-truth	"hello"	"world"	"<end>"		
Prediction 1	"hello"	"<end>"	"foo"		
Prediction 2	"hello"	"how"	"are"	"you"	"<end>"


 $loss_1$ $loss_2$ $loss_3$

Mini-batch training: padding

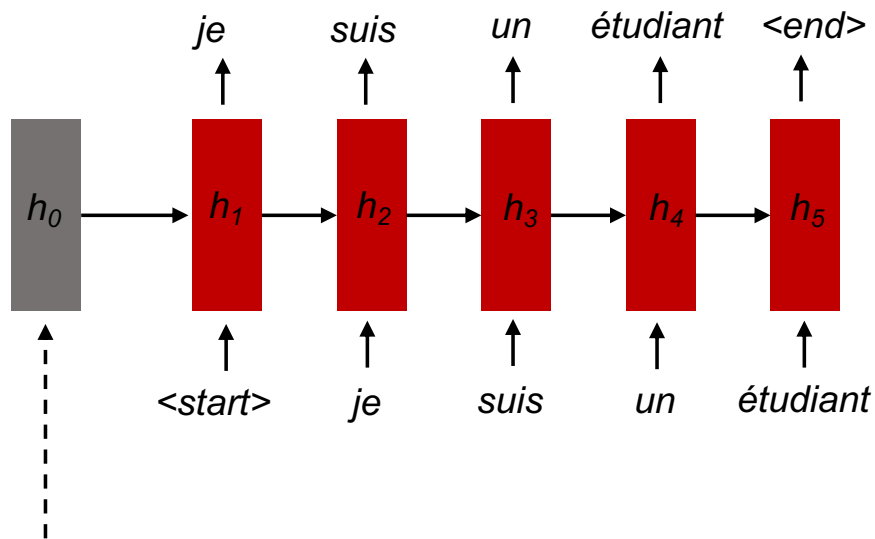
Sample 1	"hello"	"how"	"are"	"you"	"today"	"<end>"
Sample 2	"a"	"dog"	"is"	"driving"	"<end>"	"<pad>"
Sample 3	"hello"	"world"	"<end>"	"<pad>"	"<pad>"	"<pad>"

Mini-batch training: padding

Sample 1	"hello"	"how"	"are"	"you"	"today"	"<end>"
Sample 2	"a"	"dog"	"is"	"driving"	"<end>"	"<pad>"
Sample 3	"hello"	"world"	"<end>"	"<pad>"	"<pad>"	"<pad>"

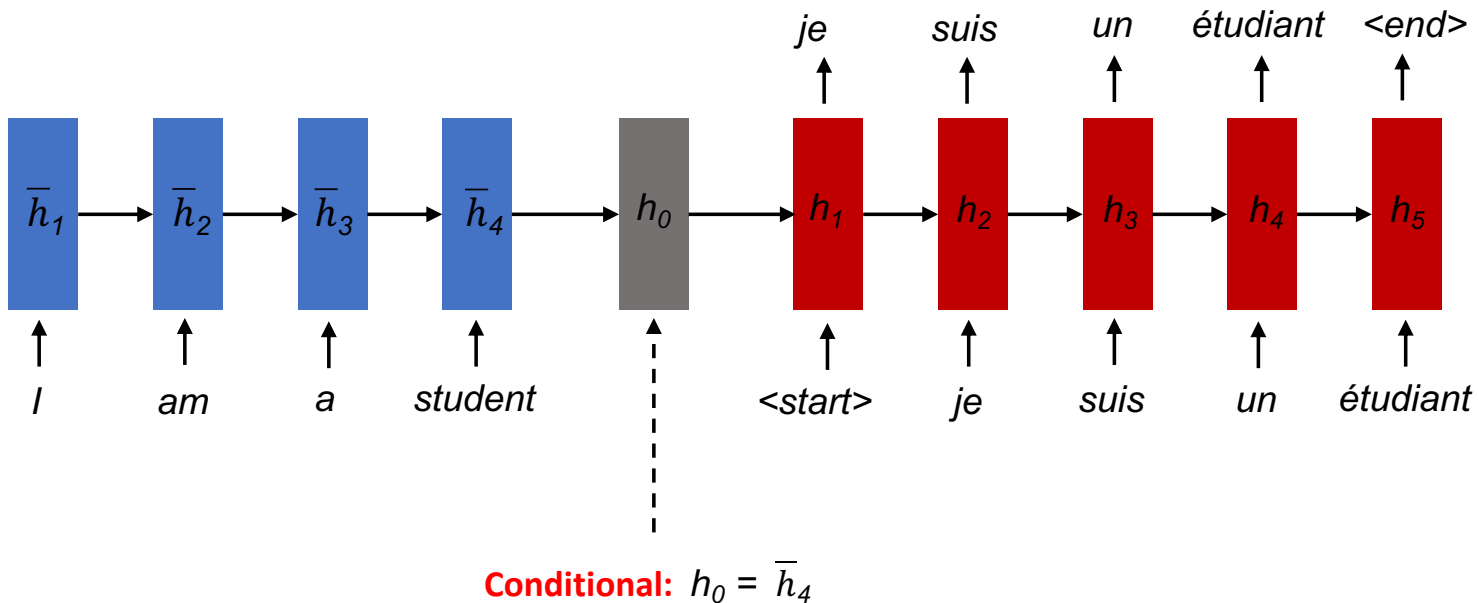
Sample 1	1.0	1.0	1.0	1.0	1.0	1.0
Sample 2	1.0	1.0	1.0	1.0	1.0	0.0
Sample 3	1.0	1.0	1.0	0.0	0.0	0.0

Generating text that makes sense: Language Model

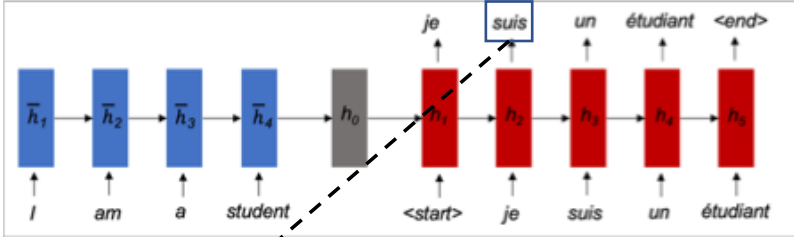


Unconditional: $h_0 = 0$

Generating text with a goal: Machine Translation



Seq2Seq model



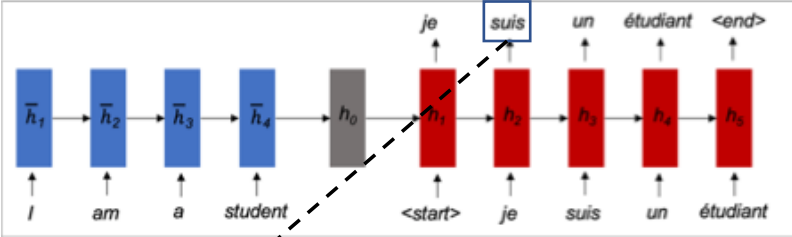
Seq2Seq:

"suis"

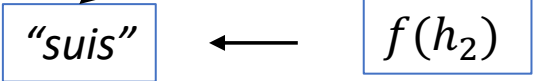


$$f(h_2) = \text{softmax}(W_S h_2)$$

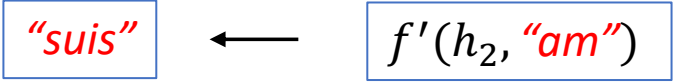
Seq2Seq model with perfect word alignments



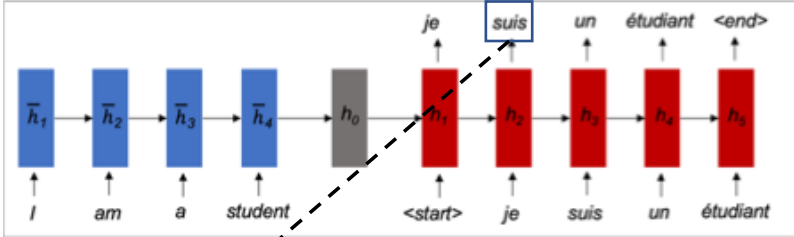
Seq2Seq:



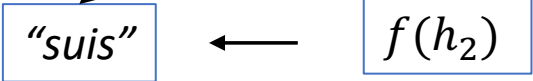
Ideally:



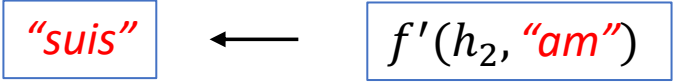
Seq2Seq model with perfect word alignments



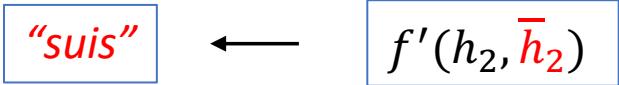
Seq2Seq:



Ideally:



Or:



Seq2Seq model with attention

Ideally: "suis" ← $f'(h_2, \bar{h}_2)$

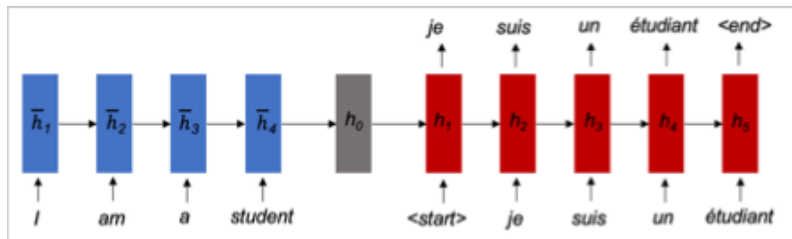
In practice: "suis" ← $f'(h_2, c_2)$

$$c_2 = \sum_{k=1}^4 w_{2,k} \bar{h}_k$$

Pray that **S**: $w_{2,2} = 1, w_{2,k \neq 2} = 0$ is true

Or train the model such that **S** is almost true

Seq2Seq model with attention



Key assumption: $\bar{h}_2 \approx h_2 \approx h_0$ – “je” $\approx \bar{h}_4$ – “je”

$$w_{2,k} = \frac{\exp(\text{score}(h_2, \bar{h}_k))}{\sum_j \exp(\text{score}(h_2, \bar{h}_j))}$$

$$\text{score}(h_2, \bar{h}_k) = h_2^T W_a \bar{h}_k$$

Seq2Seq model with attention

"suis" ← $f'(h_2, c_2) = \text{softmax}(W_s \tanh(W_c[h_2; c_2]))$

Perform much better for long sequences

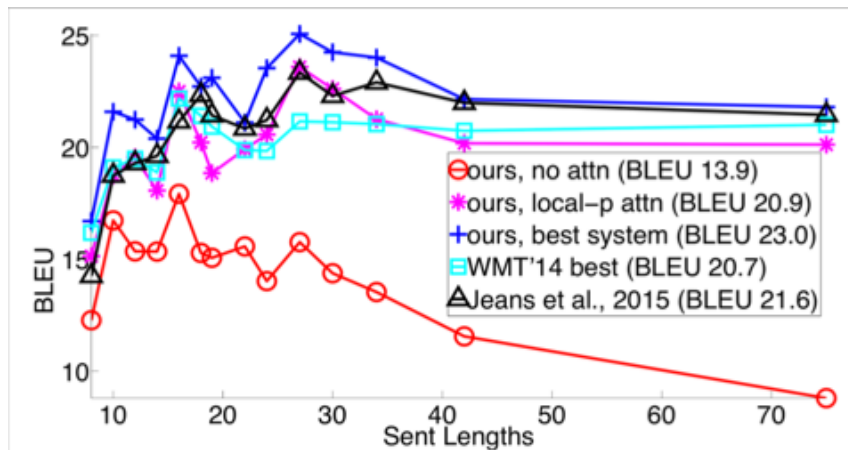
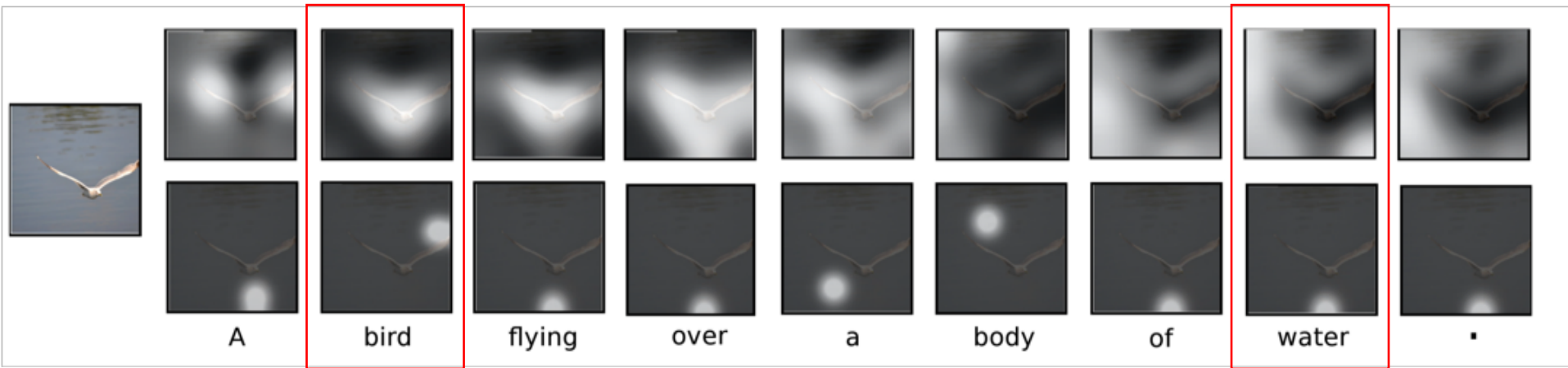


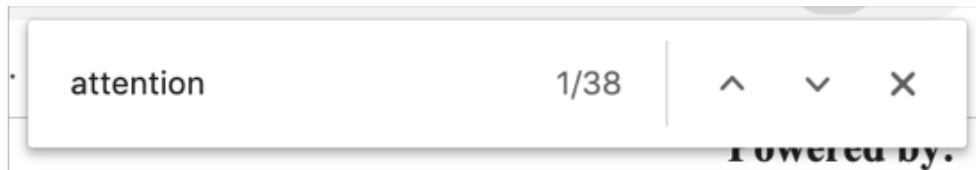
Figure 6: **Length Analysis** – translation qualities of different systems as sentences become longer.

Also very helpful in image captioning



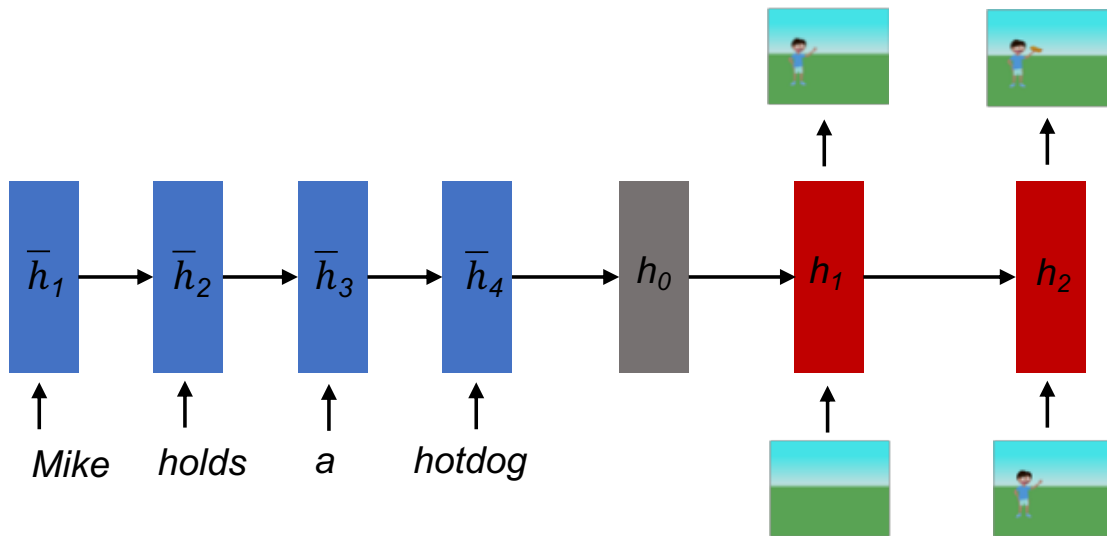
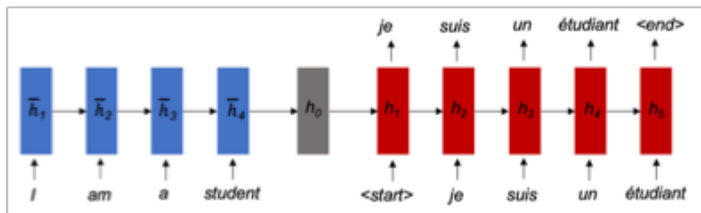
[4] Show, attend and tell: neural image caption generation with visual attention. Kelvin Xu, Jimmy Lei Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Richard Zemel, Yoshua Bengio. ICML 2015

ECCV 2018 accepted 776 papers



38 of them with “attention” in their titles

Can we do this?



Challenges

Machine Translation:

"I am a student" → "je suis un étudiant"

Text-to-Image Synthesis:

"A person is
holding a surfboard"

?

person
surfboard

person
surfboard

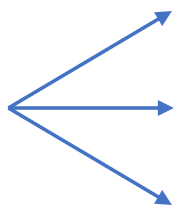
Challenges: in each step

Machine Translation:

student → *étudiant*

Text-to-Image Synthesis:

"A person is
holding a surfboard"



object category: person, surfboard

location: somewhere in the 2D world

attributes: size, pose, expression, ...

Challenges: in each step

Text-to-Image Synthesis:

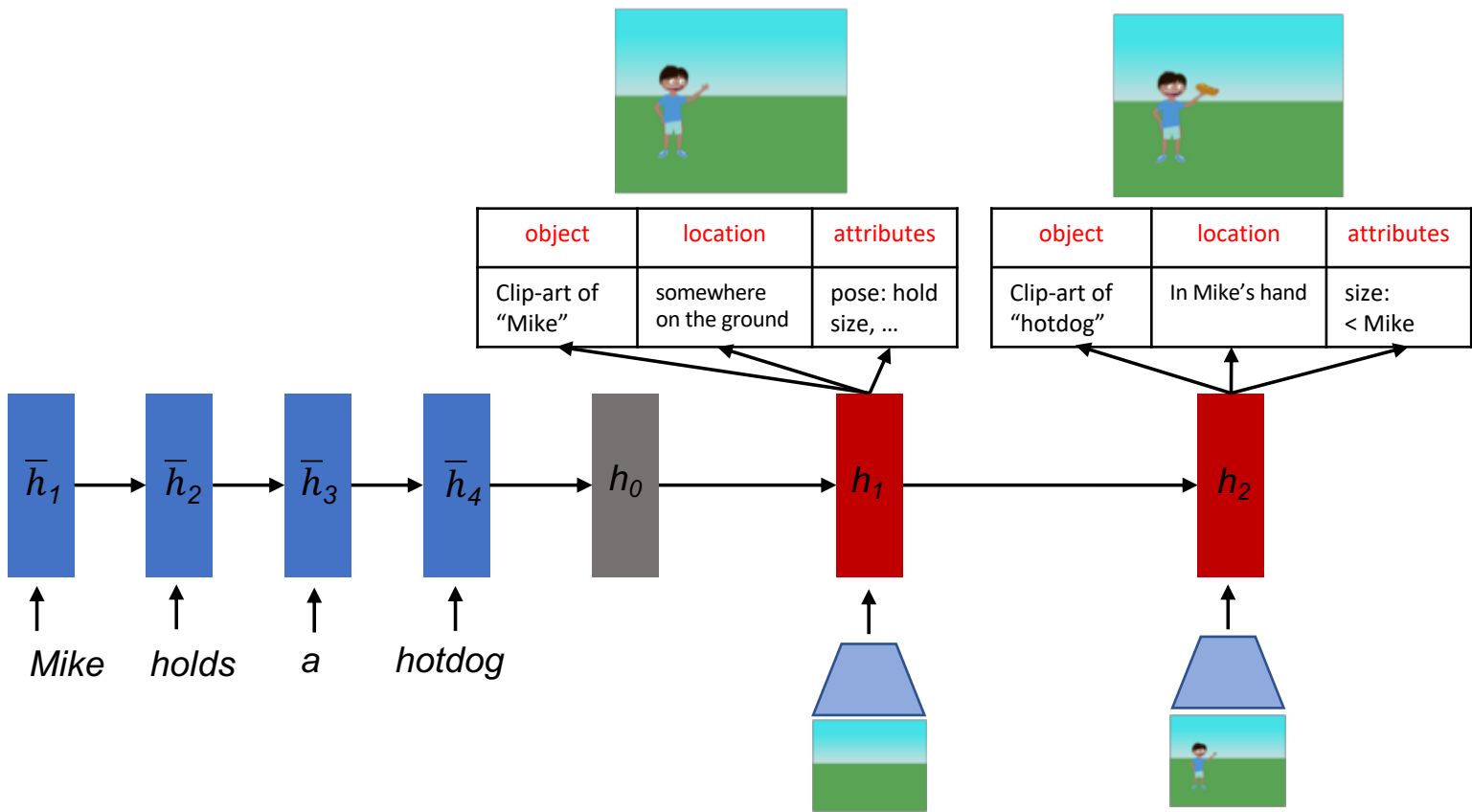
"A person is
holding a surfboard"

object category: person

location: somewhere in the 2D world

attributes: size, pose, expression, ...

Learning the distributions of
categories, locations, attributes
from the training samples



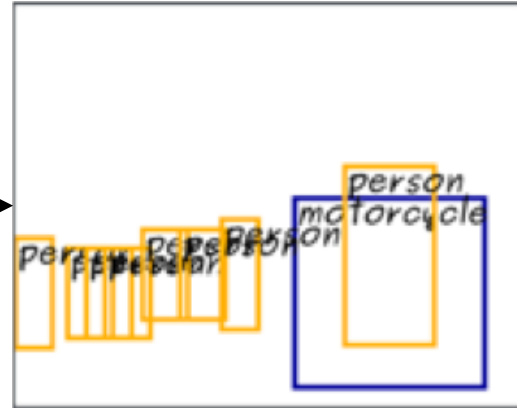
Task 1: Abstract Scene Generation

“Mike is surprised at the duck. The duck is standing on the grill. Jenny is running towards Mike and the duck.”

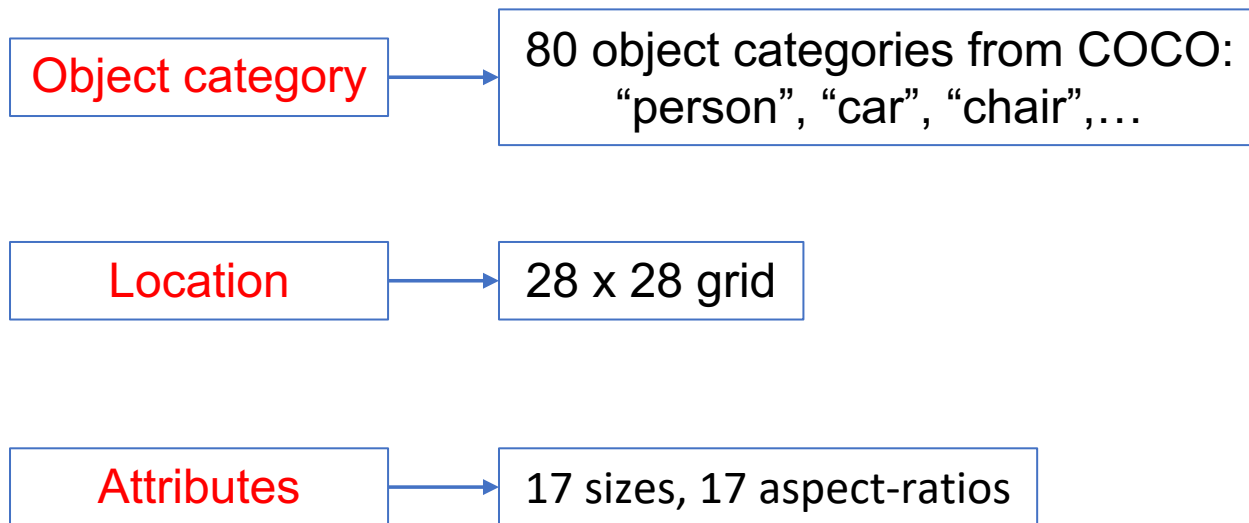


Task 2: Scene Layout Generation

“A guy on a motorcycle
with some people
watching.”



Task 2: Scene Layout Generation

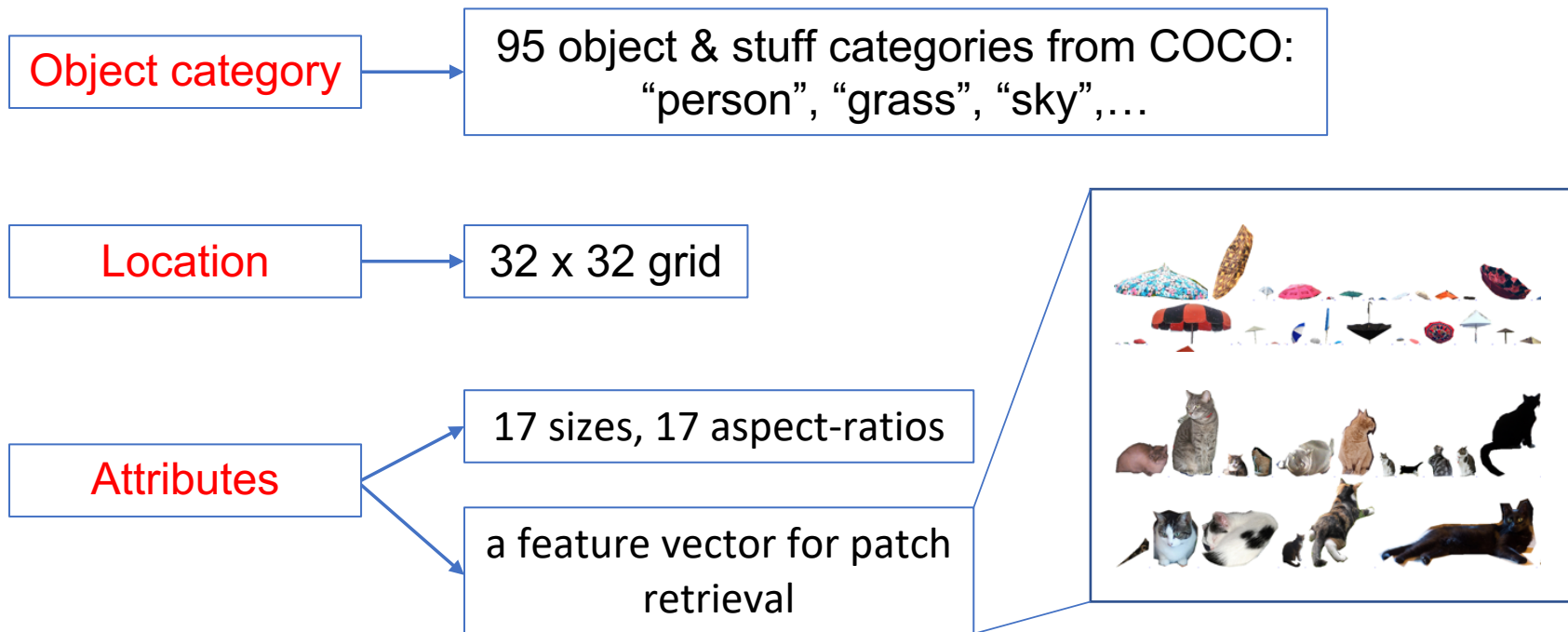


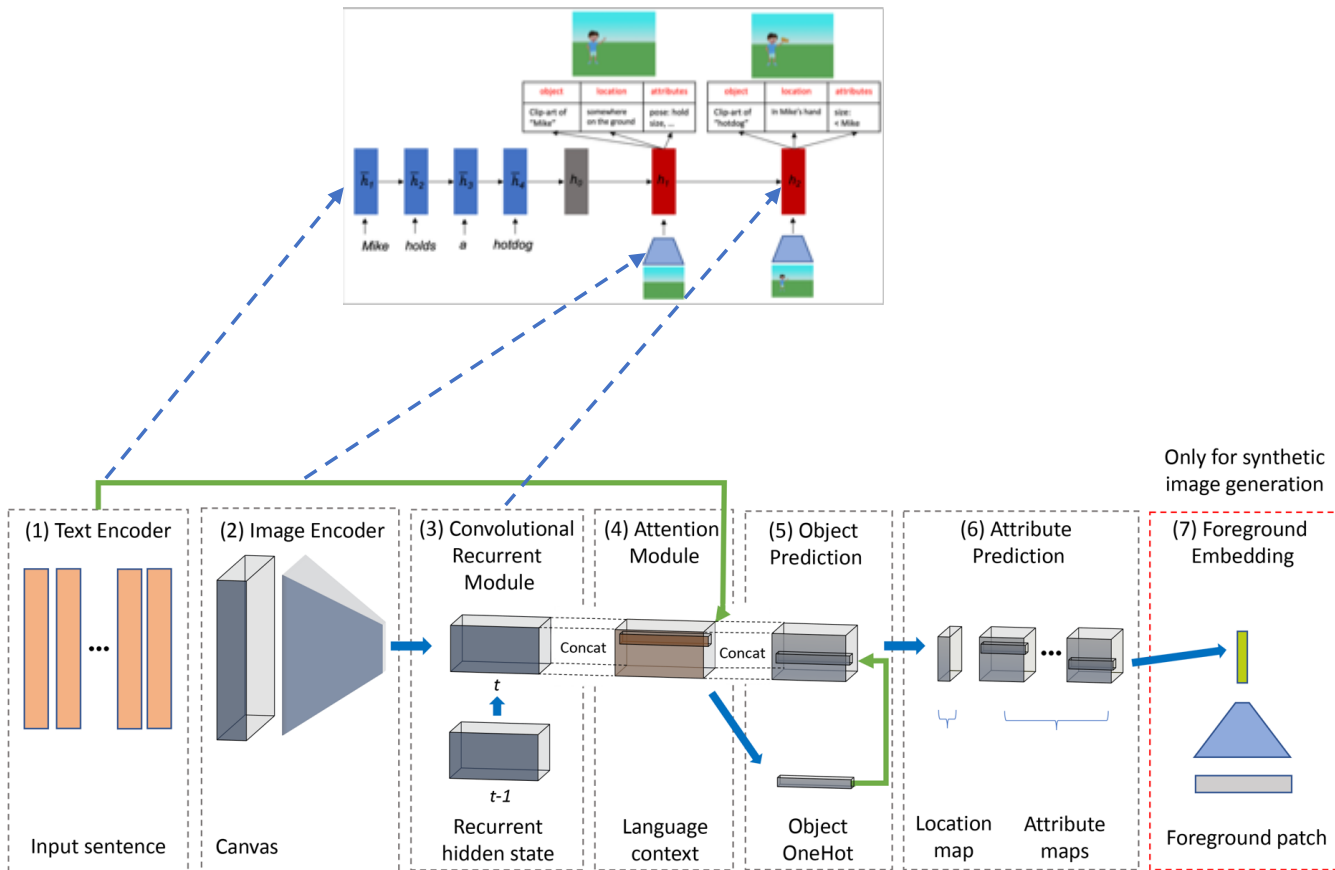
Task 3: Composite Image Generation

“Several elephants walking together in a line near water.”



Task 3: Composite Image Generation



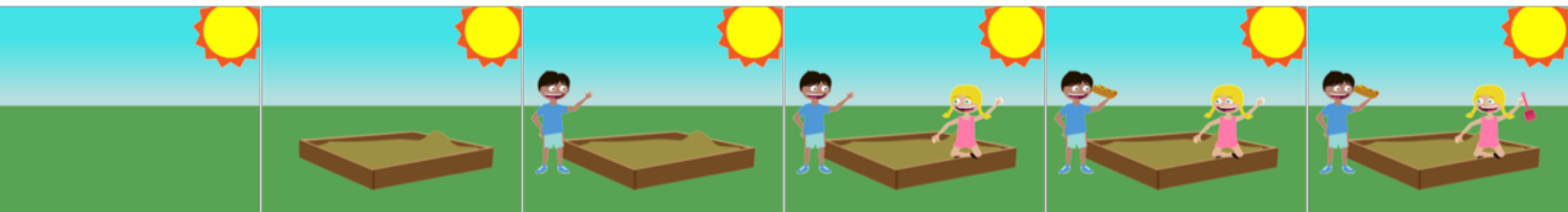


Step-by-step generation of Abstract Scene

Mike is holding a hotdog.

Jenny is sitting in the sandbox.

Jenny is holding the shovel.



object attn:
sitting sandbox holding
attribute attn:
jenny <eos> jenny

object attn:
sandbox sitting mike
attribute attn:
sandbox <eos> jenny

object attn:
mike jenny sitting
attribute attn:
holding hotdog mike

object attn:
jenny jenny mike
attribute attn:
sitting jenny holding

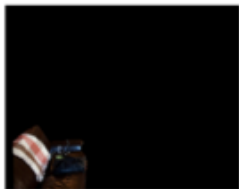
object attn:
hotdog shovel holding
attribute attn:
mike hotdog holding

object attn:
shovel holding sandbox
attribute attn:
shovel holding <eos>

Step-by-step generation of composite image

Inputs: a room with TV and some different types of couches.

couches, room, TV



room, couches, <eos>



<eos>, room, types



TV, types, different



types, of, and




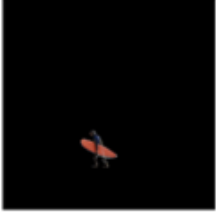




Generated sequences and the top-3 attended words in each step.




Reference image

Step-by-step generation of composite image

Inputs: a person walks on the beach, carrying a surf board.












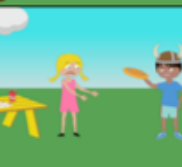




carrying, board, surf	surf, board, carrying	the, on, a
		
beach, surf, on	the, on, a	on, a, the
		

Generated sequences and the top-3 attended words in each step.



Reference image

More examples

Input	Zitnick et al. 2013	Text2Scene (w/o Attention)	Text2Scene	Reference
Jenny is wearing sunglasses. Mike is holding the red shovel. Mike is wearing a viking head.				
Mike went down the slide fast. Jenny is worried that Mike is hurt. Jenny is wearing a chef hat.				
Mike is angry at Jenny. Jenny is sad that Mike took the frisbee. The pizza is on the table.				
Jenny is holding a bucket and shovel. Mike fell off the swingset. There is rain and lightning in the sky				



More examples

Input Caption	Predicted Layout	Reference Layout	Reference Image	Input Caption	Predicted Layout	Reference Layout	Reference Image
An attractive young woman leads a grey horse through a paddock.				A couple of women ride horses through some water.			
Two giraffes in a zoo enjoy a walk and a snack.				A cat standing next to of an open refrigerator door.			
A person holding a surf board in a body of water.				This is a man riding a board in the water.			
A laptop computer a keyboard and two monitors.				A woman is riding her bike down the street in front of traffic.			
A man and a woman stand under an umbrella at a street crossing on a rainy day.				Two women walk outside, both holding up umbrellas.			

More examples

Input Caption	Real Image	SG2IM	HDGAN	AttnGAN	Text2Scene [no inpainting]	Text2Scene
A room with a TV and some different types of couches .						
A tall monitor is near a keyboard and a mouse .						
a car bridge going over a commuter train .						
Three zebras grazing in a grassy area near shrubs.						
A woman sitting on a bench with an umbrella on her head.						
A woman is riding her bike down the street in front of traffic .						

More examples

Input Caption	Text2Scene	Real Image
A woman sitting on a horse on a dirt field.		
The man is practicing his tricks on his skateboard.		
		
	Source Images	Source Images
Input Caption	Text2Scene	Real Image
A person that is eating some food on a table.		
A city bus drives down a road with other traffic.		
		
	Source Images	Source Images

Questions?